# Data-driven Characterization of Fake Content in Social Media

FAN Cheuk Pan

# Agenda

Background

Methodology

Data Acquisition

Classification Model

Application Implementation

Conclusion

# Background

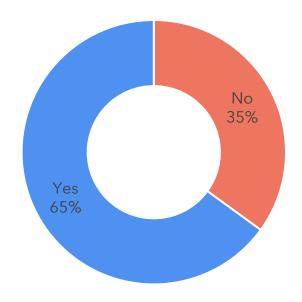Background → Problem Statement → Objectives

# Background

Fake content

Misinformation

Disinformation

# Background



CIGI Ipsos Global Survey 2019

Reported seeing fake news on social media



No 35%

Yes 65%

# Problem Statement

**Cause economic loss**
- Fake shop

**Gain advantage by bad intention**
- Gain vote in election by spread disinformation

**Cause social panic**
- False COVID-19 treatments
- Panic Buying

# Objectives

Verify fake content

Analysis of fake content

Rising awareness of fake content

# Methodology

Problem Framing → Data Acquisition → Exploratory Data Analysis

# Methodology

Fake content detector

Tool for general public

Simple Rule to Classify fake content

# System Component

**Fake Content Detector**

| Fake content detector | Simple Rule | Application |
|---|---|---|
| Feature Based Approach (Machine Learning) | Learning Based Approach (Deep Learning) | Analysis |

| | | |
|---|---|---|
| Chrome Plugin | Web Dashboard | |

# Data Acquisition – Data Source

Data Source

- FakeNewsNet
  - Gossip Cop
  - PolitiFact
- Twitter Data

# Data Acquisition – Fetch from Twitter

| Tweet |
|---|
| Content |
| Public metrics |
| Hashtag |
| URL |
| Context Domain |
| Created Time |
| Label |

| User |
|---|
| Public metrics |
| Description |
| Created time |

# Data Acquisition – Pre-Processing

Python Spacy

- Tokenization
- Remove
  - Stop word
  - Number
- Lemmatization
- Lower case

# Data Acquisition – Pre-Processing



Data Source

Fetch from data

Processed Text

# Data Acquisition – Split Data

# Exploratory Data Analysis (EDA)

| Total Record | • Label Distribution |
|---|---|
| **Sentence Length** | • Word<br>• Character |
| **Term Frequency** | • Word Cloud<br>• Word Ranking |
| **Meta data** | • Hashtag<br>• Domain (Annotated by Twitter) |

# EDA – Label ExDistribution



Label Distribution

| False (Fake) | 478069 |
|---|---|
| True (Real) | 956506 |

# EDA – Sentence Length

# EDA – Word Cloud of Content

**False Label**

**All**

**True Label**

# EDA – Word Cloud of Content

**False Label**

**All**

**True Label**

# EDA – Top 10 Hashtags



Tag Count (Top 10)

# Classification Model – Procedure



| Sentence (String) | → | Preprocessing (Spacy) | → | Frequency Embedding | → | Classification Model |
| | | | | Similarity Embedding | | |

| Tokenization | | Truncate |
| Lemmatization | | Zero-Padding |

# Word Embedding – Frequency

- Count (BOW)
- TF-IDF

# Word Embedding – Similarity

- Custom Train from Pre-trained
  - Word2Vec
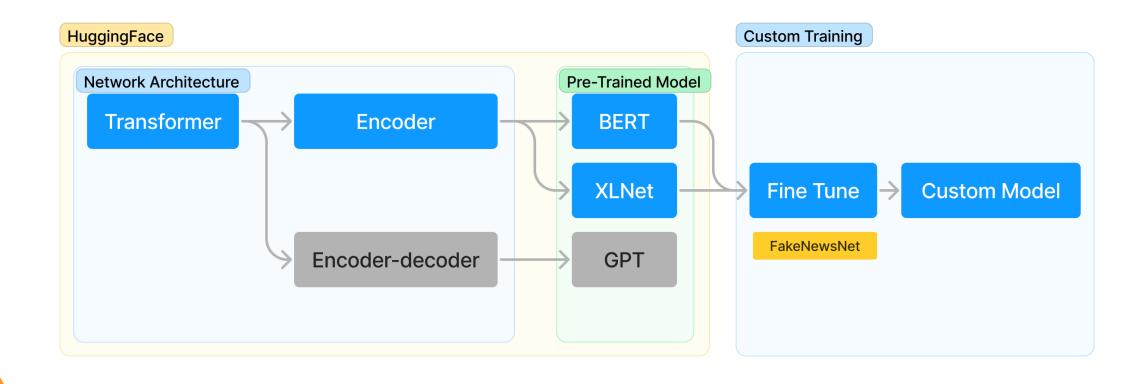    - Score: 0.450461
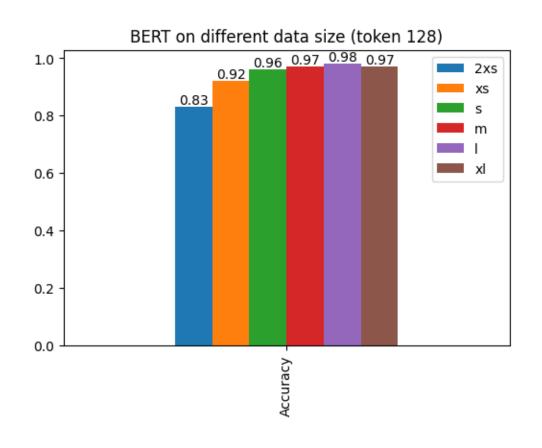  - FastText
    - Score: 0.226015

# Classification Model

# Classification Model – Result (Accuracy)

# Transformer Model

# Transformer Model – BERT on different data size



BERT on different data size (token 128)

# Transformer Model – BERT on different token size

BERT on different token size (dataset m size)

# Transformer Model – BERT and XLNet

# Model – Result conclusion

## Feature-based

- Achieve more than 80% accuracy
- Neural Network work well with Similarity Word Embedding Model

## Learning-based

- Achieve over 90% accuracy
- More data result in more
- Short-text classification, smaller token size not sightly affect the result
- Same training data size, XLNet may perform better

## Overfit

- No limitation on maximum feature

# Application

| System Architecture | → | Database | → | Chrome Plugin | → | Web Dashboard | → | Deployment |

# Recall Objectives

**Verify fake content**
- Fake content detector

**Analysis of fake content**
- Simple Rule to Classify fake content

**Rising awareness of fake content**
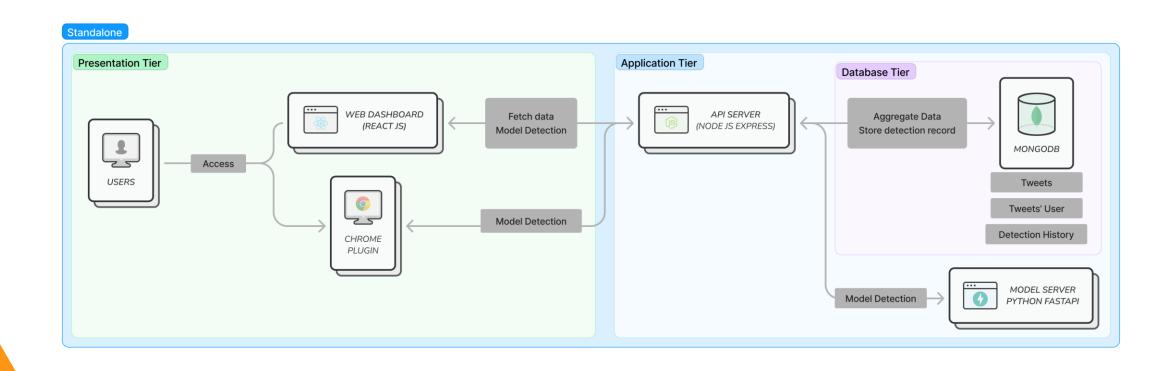- Tool for general public

# Application

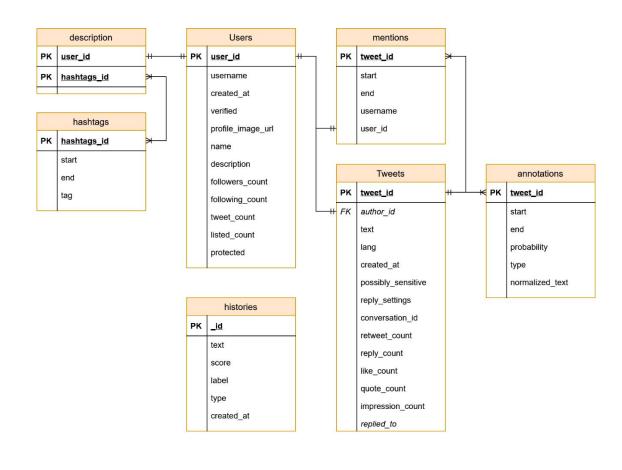## Chrome Plugin

- Perform detection

## Web Dashboard

- Perform detection
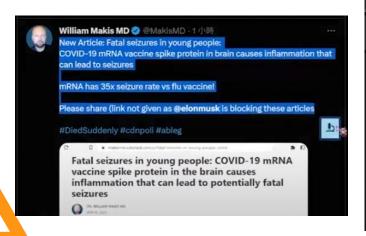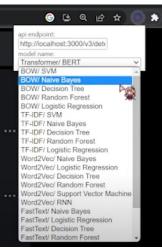- Data Visualization of dataset
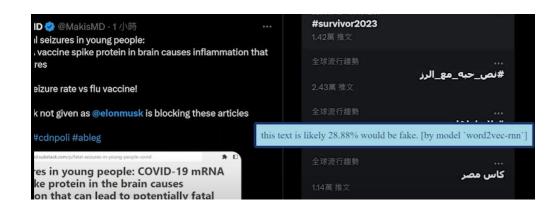
# System Architecture
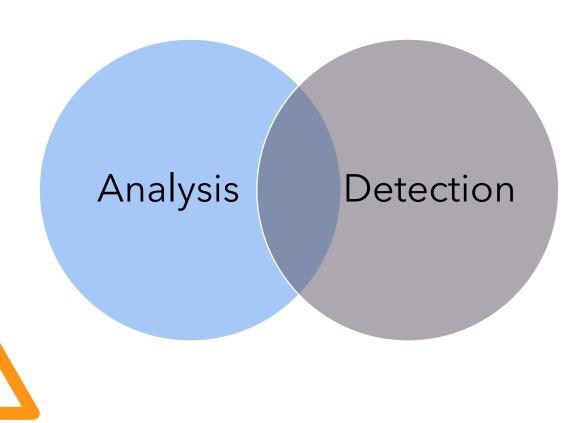
# Database

# Chrome Plugin



Select the text

Click the icon

Detect

# Web Dashboard

# Web Dashboard Detection



user | React frontend implementation | Express backend server | AWS sagemaker endpoint | Prisma client | Database

1. Input text
2. Click "Detect"
Send request with text
validate the request
Request model detection
wait for cold start of image
return detection result
store detection result
insert operation
response of operation
response of operation
return the detection result
display the result
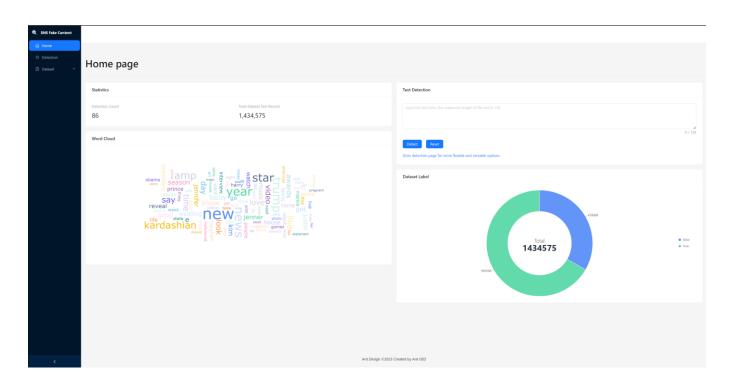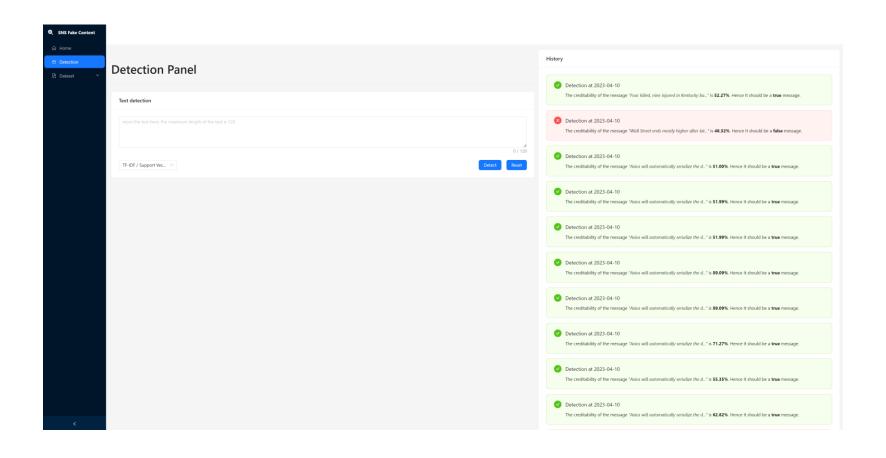
# Web Dashboard – Home Screen

- Home Page
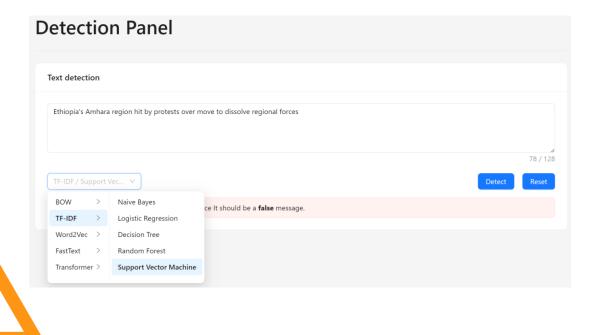    - Overview of Word Cloud
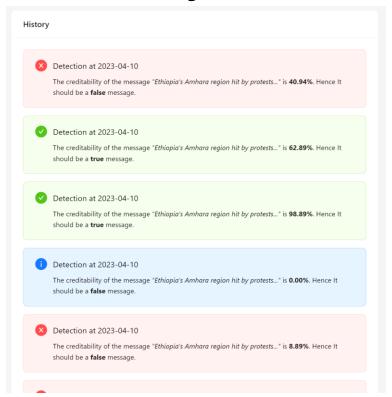    - Perform Detection

# Web Dashboard – Detection

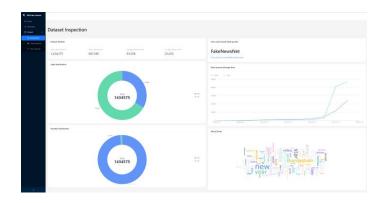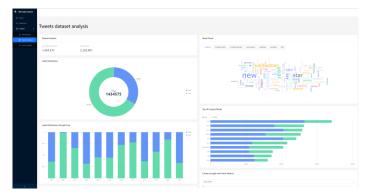# Web Dashboard – Detection

## Detection Textbox



## Detection History

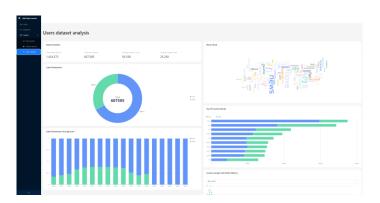# Web Dashboard – Dataset Analysis

**Dataset Information**     **Tweets' Information**     **Users' Information**
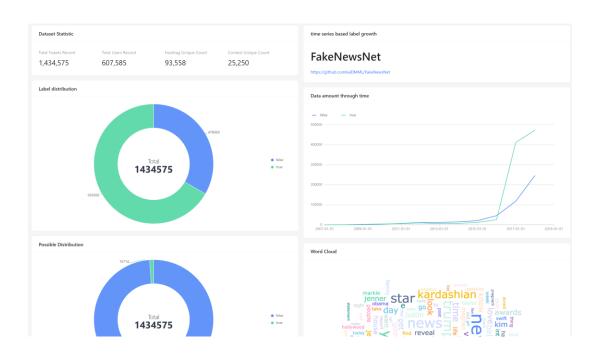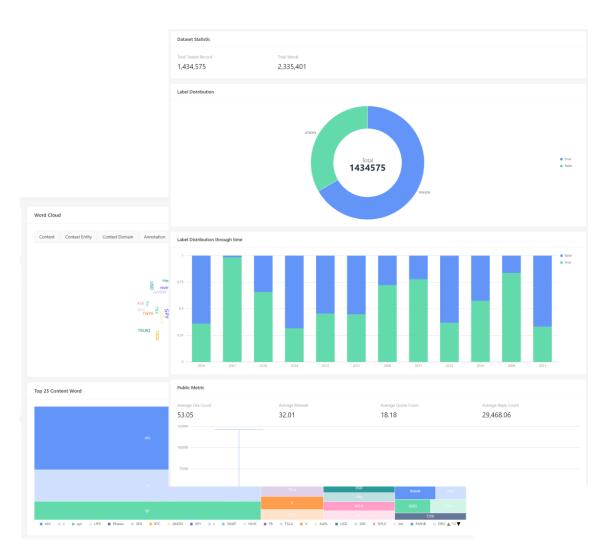
# Web Dashboard – Dataset Information

- General Dataset Information
  - Label Distribution
  - Possible Sensitive Content
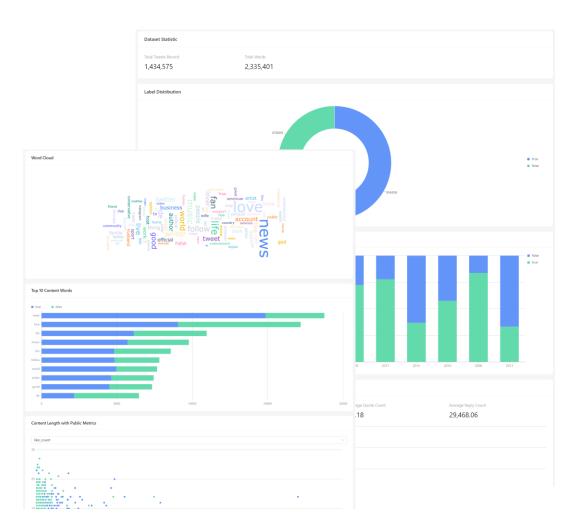  - Data Amount through time
  - Word Cloud

# Web Dashboard – Tweets Information

- Statistic
- Label Distribution
  - All
  - Through Time
- Public Metrics
- Term Frequency
  - Content
  - Entity
  - Context Domain
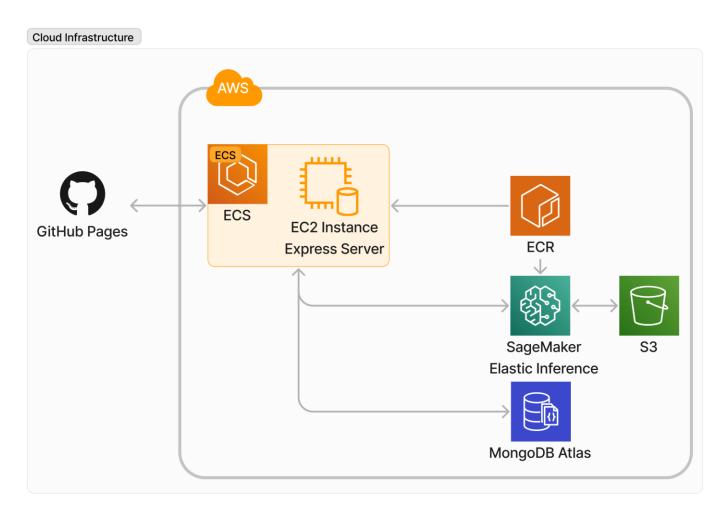  - Annotation
  - Hashtag
  - Cashtag
  - Url
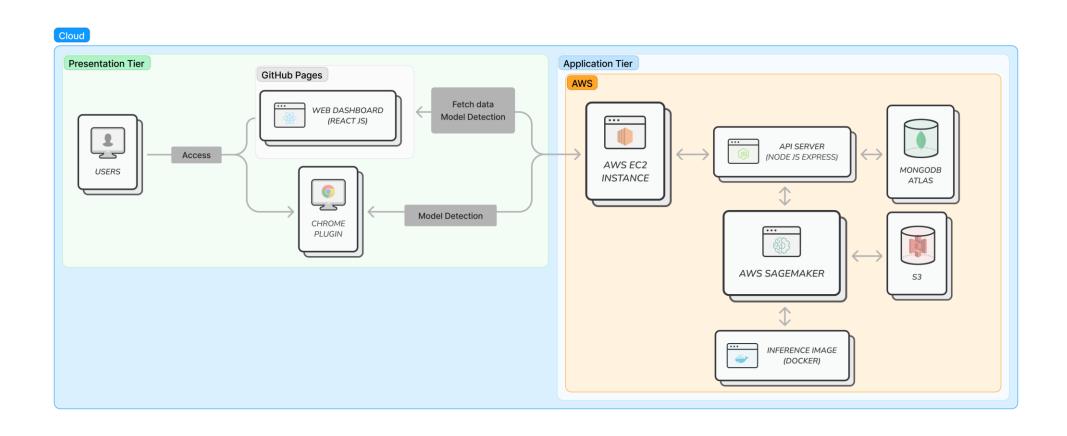
# Web Dashboard — Users Information
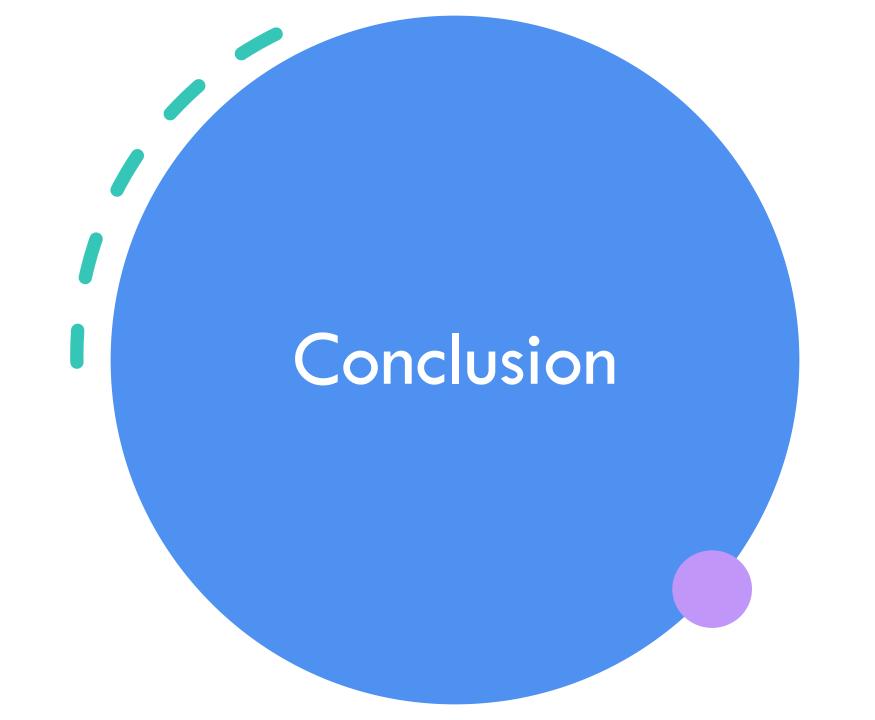
- Data Statistic
- User Description information

# Deployment

# Deployment

# Conclusion

# Future work and Improvement
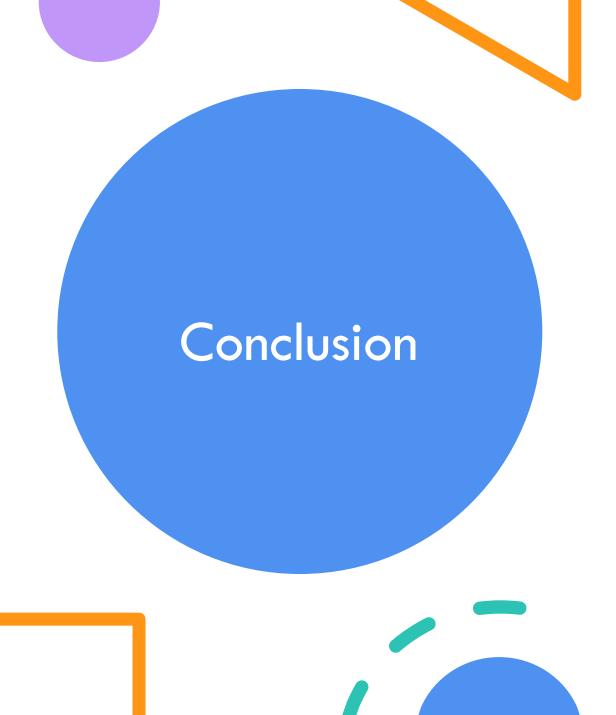
**Analysis**

- More detail on analysis
- Unsupervised learning on feature
  - Clustering

**Model**

- Expand dataset

**Application**

- Optimize detection performance
- Real time analysis
  - Scrap social media content
  - Perform detection and analysis in data pipeline

# Conclusion

## Classification Model

- Different model algorithm
- Different token size, data size

## Implementation

- Web dashboard
- Chrome Plugin

# The End

Q&A Section